

Sample Size Dependent Species Models

Mingyuan Zhou and Stephen G Walker

The University of Texas at Austin

October 14, 2014

Abstract

Motivated by the fundamental problem of measuring species diversity, this paper introduces the concept of a cluster structure to define an exchangeable cluster probability function that governs the joint distribution of a random count and its exchangeable random partitions. A cluster structure, naturally arising from a completely random measure mixed Poisson process, allows the probability distribution of the random partitions of a subset of a sample to be dependent on the sample size, a distinct and motivated feature that differs it from a partition structure. A generalized negative binomial process model is proposed to generate a cluster structure, where in the prior the number of clusters is finite and Poisson distributed, and the cluster sizes follow a truncated negative binomial distribution. We construct a nonparametric Bayesian estimator of Simpson's index of diversity under the generalized negative binomial process. We illustrate our results through the analysis of two real sequencing count datasets.

Keywords: Bayesian nonparametrics, exchangeable cluster/partition probability functions, generalized gamma process, generalized negative binomial process, generalized Chinese restaurant sampling formula, partition structure, species sampling.

M. Zhou is with the Department of Information, Risk, and Operations Management, McCombs School of Business, and S. G. Walker is with the Departments of Mathematics and Statistics & Data Sciences, the University of Texas at Austin, Austin, TX 78712, USA. *Emails:* mingyuan.zhou@mcombs.utexas.edu, s.g.walker@math.utexas.edu.

1 Introduction

A fundamental problem in biological and ecological studies is to measure the degree of diversity of a population whose individuals are classified into different groups; see Fisher et al. (1943), Simpson (1949), Hill (1973) and Magurran (2004). The rapid development of modern sequencing technologies also generates significant recent interest in the measurement of population diversity using samples summarized as the frequencies of observed sequences (Hughes et al., 2001, Shaw et al., 2008, Bunge et al., 2014, Guindani et al., 2014). The Simpson’s index of diversity, widely used to measure species evenness, is defined as the probability for two individuals randomly selected from a population to be from different groups (Simpson, 1949). Thus, if π_k denotes the population probability for an individual to be in group k , with $\sum_{k \geq 1} \pi_k = 1$, then the Simpson’s index of diversity is defined as

$$S = 1 - \sum_{k=1}^K \pi_k^2 \quad (1)$$

which is also understood to be $P(z_1 \neq z_2)$, where z_i is the group individual i is assigned to. Here, K could be finite or infinite though Simpson (1949) assumed it to be finite.

A sample estimate for (1), which is unbiased, is given by

$$\hat{S} = 1 - \sum_{k=1}^K \frac{n_k(n_k - 1)}{n(n - 1)}, \quad (2)$$

where

$$n_k = \sum_{i=1}^n \mathbf{1}(z_i = k).$$

Alongside Simpson’s index of diversity, other diversity indices have been proposed to measure species richness; see Bunge and Fitzpatrick (1993), Chao (2005) and Bunge et al. (2014) for reviews. Recent nonparametric Bayesian approaches to species diversity, focusing on the study of species richness, derive the distribution of the number of new species via n' new individuals randomly selected from the population, given a sample of size n ; see Lijoi et al. (2007a, 2008) and Favaro et al. (2009, 2013). These papers form the basis for Bayes nonparametric estimators of the Simpson’s index of diversity, as in Cerquetti (2012).

The underlying structure of the Bayesian species sampling models are built on Kingman’s concept of a partition structure, (Kingman, 1978a,b), which defines a family of consistent probability distributions for random partitions of a set $[m] := \{1, \dots, m\}$. The sampling consistency requires the probability distribution of the random partitions of a subset of size

m of a set of size $n \geq m$ to be the same for all n . More specifically, for a random partition $\Pi_m = \{A_1, \dots, A_l\}$ of the set $[m]$, where there are l clusters and each element $i \in [m]$ belongs to one and only one set A_k from Π_m , such a constraint requires that $P(\Pi_m|n) = P(\Pi_m|m)$ does not depend on n , where $P(\Pi_m|n)$ denotes the marginal partition probability for $[m]$ when it is known the sample size is n . As further developed in (Pitman, 1995, 2006), if $P(\Pi_m|m)$ depends only on the number and sizes of the (A_k) , regardless of their order, then it is called an exchangeable partition probability function (EPPF) of Π_m , expressed as $P(\Pi_m = \{A_1, \dots, A_l\}|m) = p_m(n_1, \dots, n_l)$, where $n_k = |A_k|$. The sampling consistency amounts to an addition rule (Pitman, 2006, Gneden et al., 2009) for the EPPF; that $p_1(1) = 1$ and

$$p_m(n_1, \dots, n_l) = p_{m+1}(n_1, \dots, n_l, 1) + \sum_{k=1}^l p_{m+1}(n_1, \dots, n_k + 1, \dots, n_l). \quad (3)$$

An EPPF of Π_m satisfying this constraint is considered as an EPPF of $\Pi := (\Pi_1, \Pi_2, \dots)$. For an EPPF of Π , Π_{m+1} can be constructed from Π_m by assigning element $(m+1)$ to $A_{z_{m+1}}$ based on the prediction rule as

$$z_{m+1}|\Pi_m = \begin{cases} l+1 & \text{with probability } \frac{p_{m+1}(n_1, \dots, n_l, 1)}{p_m(n_1, \dots, n_l)}, \\ k & \text{with probability } \frac{p_{m+1}(n_1, \dots, n_k+1, \dots, n_l)}{p_m(n_1, \dots, n_l)}. \end{cases}$$

A basic EPPF of Π is the Ewens sampling formula (Ewens, 1972, Antoniak, 1974). Moving beyond the Ewens sampling formula, various approaches, including the Pitman-Yor process (Perman et al., 1992, Pitman and Yor, 1997), Poisson-Kingman models (Pitman, 2003), species sampling (Pitman, 1996), stick-breaking priors (Ishwaran and James, 2001), and Gibbs-type random partitions (Gneden and Pitman, 2006), have been proposed to construct more general EPPFs of Π . See Müller and Quintana (2004), Lijoi and Prünster (2010) and Müller and Mitra (2013) for reviews. Among these approaches, there has been increasing interest in normalized random measures with independent increments (NRMIs) (Regazzini et al., 2003), where a completely random measure (Kingman, 1967, 1993) with a finite and strictly positive total random mass is normalized to construct a random probability measure. For example, the normalized gamma process is a Dirichlet process (Ferguson, 1973). More advanced completely random measures, such as the generalized gamma process of Brix (1999), can be employed to produce more general exchangeable random partitions of Π (Pitman, 2003, 2006, Lijoi et al., 2007b). However, the expressions of the EPPF and its associated prediction rule usually involve integrations that are difficult to calculate.

With respect to the Simpson's measure of diversity, it is our contention that a prior model for this quantity; i.e. $P(z_1 \neq z_2)$ should depend on n and hence we write $P(z_1 \neq z_2|n)$ meaning that in general, rather than the marginal distribution of (z_1, \dots, z_m) , with (z_{m+1}, \dots, z_n) integrated out, being independent of the sample size $n \geq m$, it actually does depend on n .

The motivation for this is that as n increases, so could the possible groups which are available for classification. It is anticipated that unknown species emerge, which is different from known species first being seen, as samples are collected. Hence, the probability, according to the experimenter's prior model, that z_1 and z_2 belong to the same group will, for example, diminish with n if, as the sample size increases, it is thought more appropriate for individuals to be reclassified into different species. In short, if all the possible species are known upfront then it is possible to classify z_1 and z_2 once and for all having seen just them. However, if there is uncertainty about the species, even whether z_1 and z_2 are the same species or not, which in life is often reality, then reassessing their classifications with n should occur and hence a model for which $P(z_1 \neq z_2)$ changes with n is motivated.

Consequently, in a Bayesian context, we will be facilitating the dependence of (z_1, \dots, z_m) , for all $m \leq n$, on n . To develop this theme, and to allow the mathematics to proceed in a neat way, and without forcing any restrictions, we make n a random object within the model.

We work at a fundamental level with a normalized completely random measure. Hence, the total (random) mass is unidentified and consequently arbitrary. We take this opportunity to use it to model the, prior to observation, random sample size n . More specifically, we model the sample size n as a Poisson random variable the mean of which is parameterized by the total random mass of a completely random measure G over a complete and separable metric space Ω . The total random mass $G(\Omega)$ is used to normalize G to obtain a random probability measure $G(\cdot)/G(\Omega)$. Linking n to $G(\Omega)$ with a Poisson distribution makes the scale of G become identifiable. With G marginalized out, the joint distribution of n and its exchangeable random partition Π_n is called an exchangeable cluster probability function (ECPF). On observing a sample of size n , we are interested in the EPPF $P(\Pi_n|n)$ and marginalizing over $n - m$ elements we would consider $P(\Pi_m|n)$. Note that distinct from a partition structure, we no longer require $P(\Pi_m|n) = P(\Pi_m|m)$ for $n > m$ in a cluster structure.

Specifically, we consider a generalized negative binomial (NB) process model where G is drawn from a generalized gamma process of Brix (1999). A draw from the generalized NB process (gNBP) represents a cluster structure with a Poisson distributed finite number of clusters, whose sizes follow a truncated NB distribution. Marginally, the sample size

follows a generalized NB distribution. These three count distributions and the prediction rule are determined by a discount, a probability and a mass parameter. These parameters are convenient to infer using the fully factorized ECPF. Since $P(\Pi_m|n) = P(\Pi_m|m)$ is often not true for $n > m$, the EPPF of the gNBP, which is derived by applying Bayes' rule on the ECPF and the generalized NB distribution, generally violates the addition rule and hence is dependent on the sample size. This EPPF will be referred as the generalized Chinese restaurant sampling formula. To generate an exchangeable random partition of $[n]$ under this EPPF, we show we could use either a Gibbs sampler or a recursively-calculated sequential prediction rule.

The layout of the paper is as follows: In Section 2 we provide all the necessary preliminary notation and a description of normalized random measures, while in Section 3 we introduce the new model for constructing sample size dependent species models. In Section 4 we apply the theory in Section 3 to the generalized negative binomial process and we present real data applications in Section 5. We end the paper with a brief conclusion and provide the proofs of theorems and corollaries in the Appendix.

2 Preliminaries

In this section we provide the mathematical foundations for an independent increment process with no Gaussian component. These are pure jump processes and for us will have finite limits so that the process can be normalized by the total sum of the jumps to provide a random distribution function. The most well known of such processes is the gamma process (see, for example, Ferguson and Klass (1972)) and we will be specifically working with a generalized gamma process in Section 2.1.

2.1 Generalized Gamma Process

The generalized gamma process, which we will denote by $\text{gGP}(G_0, a, 1/c)$, is a completely random (independent increment) measure defined on the product space $\mathbb{R}_+ \times \Omega$, where $a < 1$ is a discount parameter and $1/c$ is a scale parameter (Brix, 1999). It assigns independent infinitely divisible generalized gamma random variables $G(A_j) \sim \text{gGamma}(G_0(A_j), a, 1/c)$ to disjoint Borel sets $A_j \subset \Omega$, with Laplace transform given by

$$\mathbb{E} [e^{-\phi G(A)}] = \exp \left\{ -\frac{G_0(A)}{a} [(c + \phi)^a - c^a] \right\}. \quad (4)$$

The Lévy measure of the generalized gamma process can be expressed as

$$\nu(ds, d\omega) = \frac{1}{\Gamma(1-a)} r^{-a-1} e^{-cr} ds G_0(d\omega). \quad (5)$$

The connection between (4) and (5), not given here, is the well known form for the Laplace transform of an infinitely divisible random variable.

When $a \rightarrow 0$, we recover the gamma process, and if $a = 1/2$, we recover the inverse Gaussian process (Lijoi et al., 2005). A draw G from $\text{gGP}(G_0, a, 1/c)$ can be expressed as

$$G = \sum_{k=1}^K r_k \delta_{\omega_k},$$

with $K \sim \text{Po}(\nu^+)$ and $(r_k, \omega_k) \stackrel{i.i.d.}{\sim} \pi(ds, d\omega)$, where $r_k = G(\omega_k)$ is the weight for atom ω_k and $\pi(ds, d\omega)\nu^+ \equiv \nu(ds, d\omega)$. Except where otherwise specified, we only consider $a < 1$ and $c > 0$. If $0 \leq a < 1$, since the Poisson intensity $\nu^+ = \nu(\mathbb{R}_+ \times \Omega) = \infty$ (i.e., $K = \infty$ a.s.) and $\int_{\mathbb{R}_+ \times \Omega} \min\{1, s\} \nu(ds, d\omega)$ is finite, a draw from $\text{gGP}(G_0, a, 1/c)$ consists of countably infinite atoms. On the other hand, if $a < 0$, then $\nu^+ = -\gamma_0 c^a/a$ and thus $K \sim \text{Po}(-\gamma_0 c^a/a)$ (i.e., K is finite a.s.) and $r_k \stackrel{i.i.d.}{\sim} \text{Gamma}(-a, 1/c)$. This process will be seen again in Section 4.

2.2 Normalized Random Measures

A NRMI model (Regazzini et al., 2003) is a normalized completely random measure

$$\tilde{G} = G/G(\Omega)$$

where $G(\Omega) = \sum_{k=1}^K r_k$ is the total random mass, which is required to be finite and strictly positive. Note that the strict positivity of $G(\Omega)$ implies that $\nu^+ = \infty$ and hence $K = \infty$ a.s. (Regazzini et al., 2003, Lijoi and Prünster, 2010). For us we will not necessarily be assuming that $K = \infty$ a.s. In fact our model is such that $K = 0 \iff n = 0$, which is coherent, and, moreover, $P(K = 0 | n > 0) = 0$.

Here we describe how the random allocations of individuals to groups are distributed based on the independent random jumps of the generalized gamma process. With a random draw $G = \sum_{k=1}^K r_k \delta_{\omega_k}$, by introducing a categorical latent variable z with $P(z = k | G) = r_k/G(\Omega)$, when a sample of size n is observed we have

$$p(\mathbf{z} | G, n) = \prod_{i=1}^n \frac{r_{z_i}}{\sum_{k=1}^K r_k} = \left(\sum_{k=1}^K r_k \right)^{-n} \prod_{k=1}^K r_k^{n_k}, \quad (6)$$

where $\mathbf{z} = (z_1, \dots, z_n)$ is a sequence of categorical random variables indicating the cluster

memberships, $n_k = \sum_{i=1}^n \mathbf{1}(z_i = k)$ is the number of data points assigned to category k , and $n = \sum_{k=1}^K n_k$. A random partition Π_n of $[n]$ is defined by the ties between the (z_i) . So at this point, (6) is standard.

Now (6) exhibits a lack of identifiability in that the scale of the (r_k) is arbitrary; the model is the same if we set $\tilde{r}_k = \kappa r_k$ for any $\kappa > 0$. Hence, the total mass $\sum_{k=1}^K r_k$ is unidentified.

Additionally, for reasons outlined in Section 1, we want, having marginalized out G , for n to remain, and for us to have $p(\mathbf{z}|n)$ to remain. For the standard models, when G is integrated out, n disappears and we have $p(\mathbf{z})$ depending solely on the parameters of the model.

We solve both these issues by allowing n to depend on G via

$$p(n|G) = \text{Po}[G(\Omega)],$$

from which we have independently

$$p(n_k|G) = \text{Po}(r_k).$$

We note here then that the prior model is for $p(n, G)$ and, consequently, $p(G|n)$ means G depends on n ; i.e. for each n we will have a different random measure for G .

We provide in Section 3 the general form for the prior $p(\mathbf{z}|n)$ and in Section 4 the specific case when G is a generalized gamma process. In Section 5 we use MCMC methods to estimate the posterior values of Simpson's index of diversity using real sequence frequency count data.

Posterior inference via MCMC is also simplified by our approach. Following James et al. (2009), a specific auxiliary variable $T > 0$, with $p_T(t|n, G(\Omega)) = \text{Gamma}(n, 1/G(\Omega))$, can be introduced to yield a fully factorized likelihood, stimulating the development of a number of posterior simulation algorithms including Griffin and Walker (2011), Barrios et al. (2012) and Favaro and Teh (2013). Marginalizing out G and then T from that fully factorized likelihood leads to an EPPF of Π (Pitman, 2003, 2006, Lijoi et al., 2007b). However, the prediction rule of the EPPF may not be easy to calculate.

3 Structure of Model

As has been previously mentioned, we link the sample size n to the total random mass of G with a Poisson distribution;

$$p(n|G) = \text{Po}[G(\Omega)]. \tag{7}$$

Three distinct constructions have the same joint distribution of the total number of customers and their exchangeable random partition:

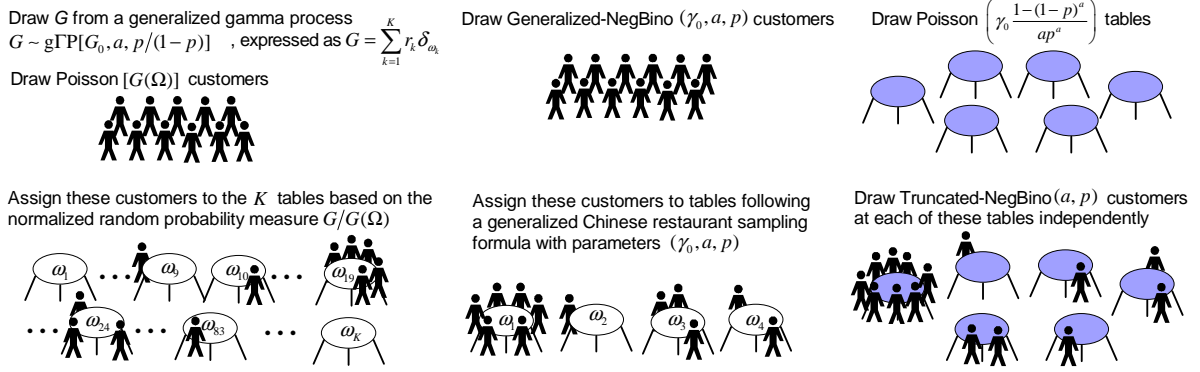


Figure 1: The cluster structure of the generalized NB process can be either constructed by assigning $\text{Pois}[G(\Omega)]$ number of customers to tables following a normalized generalized gamma process $G/G(\Omega)$, where $G \sim \text{gGP}[G_0, a, p/(1-p)]$, or constructed by assigning $n \sim \text{gNB}(\gamma_0, a, p)$ number of customers to tables following a generalized Chinese restaurant sampling formula $z \sim \text{gCRSF}(n, \gamma_0, a, p)$, where $\gamma_0 = G_0(\Omega)$. A equivalent cluster structure can also be generated by first drawing $\text{Pois}(\gamma_0 \frac{1-(1-p)^a}{ap^a})$ number of tables, and then drawing $\text{TNB}(a, p)$ number of customers independently at each table.

Since the n data points are clustered according to the normalized random probability measure $G/G(\Omega)$, we have the equivalent sampling mechanism given by

$$p(n_k|G) = \text{Po}(r_k) \quad \text{independently for } k = 1, 2, \dots,$$

and, since $n = \sum_k n_k$, we obviously recover (7).

Therefore, we link directly the cluster sizes (n_k) to the weights (r_k) with independent Poisson distributions, which is in itself an appealing intuitive feature. The mechanism to generate a sample of arbitrary size is now well defined and G is no longer scaled freely. The new construction also allows $G(\Omega) = 0$, for which $n \equiv 0$ a.s. Allowing $G(\Omega) = 0$ with a nonzero probability relaxes the requirement of $\nu^+ = \infty$ (i.e., $K = \infty$ a.s.).

A key insight of this paper is that a completely random measure mixed Poisson process produces a cluster structure that is identical in distribution to both (i) the one produced by assigning the total random count of the Poisson process into exchangeable random partitions, using the random probability measure normalized from that completely random measure, and (ii) the one produced by assigning the total (marginal) random count n of the mixed Poisson process into exchangeable random partitions using an EPPF of Π_n . For example, when the generalized gamma process (Brix, 1999) is used as the completely random measure in this setting, our key discoveries are summarized in Figure 1, which will be discussed further in Section 4.

We note that Zhou et al. (2012) and Zhou and Carin (2013) have explored related ideas

to mix a gamma or beta process with a negative binomial process, and use that hierarchical process for mixture modeling of grouped data. Yet the authors marginalized neither the beta nor gamma process due to technical difficulties and relied on finite truncation for inference. We will discuss at the end of the paper that the ideas and techniques developed in this paper serve as the foundation for the authors to develop priors for random count matrices and understand the marginal combinatorial structures of the beta-negative binomial process.

In the following theorem, we establish the marginal model for the (n_k) with G marginalized out. The proof for this theorem is provided in the Appendix.

Theorem 1 (Compound Poisson Process). *It is that the G mixed Poisson process is also a compound Poisson process; a random draw of which can be expressed as*

$$X(\cdot) = \sum_{k=1}^l n_k \delta_{\omega_k}(\cdot) \quad \text{with } l \sim \text{Po} \left[G_0(\Omega) \int_0^\infty (1 - e^{-s}) \rho(ds) \right],$$

and independently

$$P(n_k = j) = \frac{\int_0^\infty s^j e^{-s} \rho(ds)}{j! \int_0^\infty (1 - e^{-s}) \rho(ds)} \quad \text{for } j = 1, 2, \dots$$

where $\int_0^\infty (1 - e^{-s}) \rho(ds) < \infty$ is a condition required for the characteristic functions of G to be well defined, $\omega_k \stackrel{iid}{\sim} g_0$ and $g_0(d\omega) = G_0(d\omega)/G_0(\Omega)$.

The compound Poisson representation dictates the model to have a Poisson distributed finite number of clusters, whose sizes follow a positive discrete distribution. The mass parameter $\gamma_0 = G_0(\Omega)$ has a linear relationship with the expected number of clusters, but has no direct impact on the cluster-size distribution. Note that a draw from G contains $K < \infty$ or $K = \infty$ atoms a.s., but only l of them would be associated with nonzero counts if G is mixed with a Poisson process. Since the cluster indices are unordered and exchangeable, without loss of generality, in the following discussion, we relabel the atoms with nonzero counts in order of appearance from 1 to l and then $z_i \in \{1, \dots, l\}$ for $i = 1, \dots, n$, with $n_k > 0$ if and only if $1 \leq k \leq l$ and $n_k = 0$ if $k > l$.

Corollary 2 (Exchangeable Cluster/Partition Probability Functions). *The model has a fully factorized exchangeable cluster probability function (ECPF) as*

$$p(\mathbf{z}, n | \gamma_0, \rho) = \frac{\gamma_0^l}{n!} \exp \left\{ \gamma_0 \int_0^\infty (e^{-s} - 1) \rho(ds) \right\} \prod_{k=1}^l \int_0^\infty s^{n_k} e^{-s} \rho(ds),$$

the marginal distribution for the sample size $n = X(\Omega)$ has probability generating function

$$\mathbb{E}[t^n | \gamma_0, \rho] = \exp \left\{ \gamma_0 \int_0^\infty (e^{-(1-t)s} - 1) \rho(ds) \right\}$$

and probability mass function

$$p_N(n | \gamma_0, \rho) = \frac{d^n(\mathbb{E}[t^n | \gamma_0, \rho])}{dt^n} \Big|_{t=0},$$

and an exchangeable partition probability function (EPPF) of Π_n as

$$p(\mathbf{z} | n, \gamma_0, \rho) = p(\mathbf{z}, n | \gamma_0, \rho) / p_N(n | \gamma_0, \rho).$$

The proof of this is straightforward given the representation in Theorem 1 and given the one-to-many-mapping combinatorial coefficient taking (n_1, \dots, n_l, l) to (z_1, \dots, z_n, n) is

$$\frac{l!}{n!} \prod_{k=1}^l n_k!.$$

Corollary 3 (Prediction Rule). *Let l^{-i} represent the number of clusters in $\mathbf{z}^{-i} := \mathbf{z} \setminus z_i$ and $n_k^{-i} := \sum_{j \neq i} \mathbf{1}(z_j = k)$. We can express the prediction rule of the model as*

$$P(z_i = k | \mathbf{z}^{-i}, n, \gamma_0, \rho) \propto \begin{cases} \frac{\int_0^\infty s^{n_k^{-i}+1} e^{-s} \rho(ds)}{\int_0^\infty s^{n_k^{-i}} e^{-s} \rho(ds)}, & \text{for } k = 1, \dots, l^{-i}; \\ \gamma_0 \int_0^\infty s e^{-s} \rho(ds), & \text{if } k = l^{-i} + 1. \end{cases}$$

This prediction rule can be used to simulate an exchangeable random partition of $[n]$ via Gibbs sampling.

The proof for this Corollary is provided in the Appendix. In the next section we will study a particular process: the generalized negative binomial process, whose ECPF has a simple analytic expression and whose exchangeable random partitions can not only be simulated via Gibbs sampling using the above prediction rule, but also be sequentially constructed using a recursively calculated prediction rule.

4 Generalized Negative Binomial Process

In the following discussion, we study the generalized NB process (gNBP) model where $G \sim \text{gGP}[G_0, a, p/(1-p)]$ with $a < 0$, $a = 0$ or $0 < a < 1$. Here we apply the results in Section 3

to this specific case. Using (5), we have

$$\int_0^\infty s^n e^{-s} \rho(ds) = \frac{\Gamma(n-a)}{\Gamma(1-a)} p^{n-a} \quad \text{and} \quad \int_0^\infty (1-e^{-s}) \rho(ds) = \frac{1-(1-p)^a}{ap^a}.$$

Marginalizing out λ from $n|\lambda \sim \text{Po}(\lambda)$ with $\lambda \sim \text{gGamma}[\gamma_0, a, p/(1-p)]$, leads to a generalized NB distribution; i.e. $n \sim \text{gNB}(\gamma_0, a, p)$, with shape parameter γ_0 , discount parameter $a < 1$, and probability parameter p . Denote by \sum_* as the summation over all sets of positive integers (n_1, \dots, n_l) with $\sum_{k=1}^l n_k = n$. As derived in the Appendix, the probability mass function (PMF) of the generalized NB distribution can be expressed as

$$p_N(n|\gamma_0, a, p) = \frac{p^n}{n!} e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}} \sum_{l=0}^n \gamma_0^l p^{-al} S_a(n, l), \quad (8)$$

where

$$S_a(n, l) = \frac{n!}{l!} \sum_* \prod_{k=1}^l \frac{\Gamma(n_k - a)}{n_k! \Gamma(1-a)} = \frac{1}{l! a^l} \sum_{k=0}^l (-1)^k \binom{l}{k} \frac{\Gamma(n - ak)}{\Gamma(-ak)} \quad (9)$$

are generalized Stirling numbers of the first kind (Charalambides, 2005, Pitman, 2006), which can be recursively calculated via $S_a(n, 1) = \Gamma(n-a)/\Gamma(1-a)$, $S_a(n, n) = 1$ and $S_a(n+1, l) = (n-al)S_a(n, l) + S_a(n, l-1)$. Note that when $-ak$ is a nonnegative integer, $\Gamma(-ak)$ is not well defined but $\Gamma(n-ak)/\Gamma(-ak) = \prod_{i=0}^{n-1} (i-ak)$ is still well defined.

Marginalizing out G in the generalized gamma process mixed Poisson process

$$X|G \sim \text{PP}(G) \quad \text{and} \quad G \sim \text{gGP}[G_0, a, p/(1-p)] \quad (10)$$

leads to a generalized NB process $X \sim \text{gNBP}(G_0, a, p)$, such that for each $A \subset \Omega$, $X(A) \sim \text{gNB}(G_0(A), a, p)$. This process is also a compound Poisson process as

$$X(\cdot) = \sum_{k=1}^l n_k \delta_{\omega_k}(\cdot), \quad l \sim \text{Po}\left(\gamma_0 \frac{1-(1-p)^a}{ap^a}\right), \quad n_k \stackrel{iid}{\sim} \text{TNB}(a, p), \quad \omega_k \stackrel{iid}{\sim} g_0,$$

where $\text{TNB}(a, p)$ denotes a truncated NB distribution, with PMF

$$p_U(u|a, p) = \frac{\Gamma(u-a)}{u! \Gamma(-a)} \frac{p^u (1-p)^{-a}}{1-(1-p)^{-a}}, \quad u = 1, 2, \dots \quad (11)$$

The ECPF of the gNBP model is given by

$$p(\mathbf{z}, n | \gamma_0, a, p) = \frac{1}{n!} e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}} \gamma_0^l p^{n-al} \prod_{k=1}^l \frac{\Gamma(n_k - a)}{\Gamma(1 - a)}. \quad (12)$$

The EPPF of Π_n is the ECPF in (12) divided by the marginal distribution of n in (8), given by

$$p(\mathbf{z} | n, \gamma_0, a, p) = \frac{\gamma_0^l p^{-al}}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)} \prod_{k=1}^l \frac{\Gamma(n_k - a)}{\Gamma(1 - a)}. \quad (13)$$

We define the EPPF in (13) as the generalized Chinese restaurant sampling formula (gCRSF), and we denote a random draw under this EPPF as

$$\mathbf{z} | n \sim \text{gCRSF}(n, \gamma_0, a, p).$$

The conditional distribution of the cluster number in a sample of size n can be expressed as

$$p_L(l | n, \gamma_0, a, p) = \frac{1}{l!} \sum_{*} \frac{n!}{\prod_{k=1}^l n_k} p(\mathbf{z} | n, \gamma_0, a, p) = \frac{\gamma_0^l p^{-al} S_a(n, l)}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)}. \quad (14)$$

Note that if $a \rightarrow 0$, we recover, from (13), the Ewens sampling formula which is the EPPF of the Chinese restaurant process (CRP) (Aldous, 1983). The prediction rule for the EPPF in (13) can be expressed as

$$P(z_i = k | \mathbf{z}^{-i}, n, \gamma_0, a, p) \propto \begin{cases} n_k^{-i} - a, & \text{for } k = 1, \dots, l^{-i}; \\ \gamma_0 p^{-a}, & \text{if } k = l^{-i} + 1. \end{cases} \quad (15)$$

This prediction rule can be used in a Gibbs sampler to simulate an exchangeable random partition $\mathbf{z} | n \sim \text{gCRSF}(n, \gamma_0, a, p)$ of $[n]$. However, a large number of Gibbs sampling iterations may be required to generate an unbiased sample from this EPPF. Below we present a sequential construction for this EPPF.

Marginalizing out z_n from (13), we have

$$\begin{aligned} p(z_{1:n-1} | n, \gamma_0, a, p) &= p(z_{1:n-1} | n-1, \gamma_0, a, p) \\ &\times \frac{\sum_{\ell=0}^{n-1} \gamma_0^\ell p^{-a\ell} S_a(n-1, \ell)}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)} [\gamma_0 p^{-a} + (n-1) - a l_{(n-1)}], \end{aligned}$$

where $z_{1:i} := \{z_1, \dots, z_i\}$, $l_{(i)}$ denotes the number of partitions in $z_{1:i}$ and $l_{(n)} = l$. Further

marginalizing out z_{n-1}, \dots, z_{i+1} , we have

$$\begin{aligned} p(z_{1:i}|n, \gamma_0, a, p) &= p(z_{1:i}|i, \gamma_0, a, p) \frac{\sum_{\ell=0}^i \gamma_0^\ell p^{-a\ell} S_a(i, \ell)}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)} R_{n, \gamma_0, a, p}(i, l_{(i)}) \\ &= \frac{R_{n, \gamma_0, a, p}(i, l_{(i)}) \gamma_0^{l_{(i)}} p^{-al_{(i)}}}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)} \prod_{k: n_{k, (i)} > 0} \frac{\Gamma(n_{k, (i)} - a)}{\Gamma(1 - a)}, \end{aligned} \quad (16)$$

where $n_{k, (i)} := \sum_{j=1}^i \mathbf{1}(z_j = k)$; $R_{n, \gamma_0, a, p}(i, j) \equiv 1$ if $i = n$ and is recursively calculated for $i = n - 1, n - 2, \dots, 1$ with

$$R_{n, \gamma_0, a, p}(i, j) = R_{n, \gamma_0, a, p}(i + 1, j)(i - aj) + R_{n, \gamma_0, a, p}(i + 1, j + 1)\gamma_0 p^{-a}. \quad (17)$$

We name (16) as a size-dependent EPPF as its distribution on an exchangeable random partition of $[i]$ is a function of the sample size n . Note that if $a = 0$, then

$$\frac{\sum_{l=0}^i \gamma_0^l p^{-al} S_a(i, l)}{\sum_{l=0}^n \gamma_0^l p^{-al} S_a(n, l)} = \frac{\sum_{l=0}^i \gamma_0^l |s(i, l)|}{\sum_{l=0}^n \gamma_0^l |s(n, l)|} = \frac{\Gamma(i + \gamma_0)}{\Gamma(n + \gamma_0)}$$

and $R_{n, \gamma_0, a=0, p}(i, l) = \frac{\Gamma(n + \gamma_0)}{\Gamma(i + \gamma_0)}$, and hence $p(z_{1:i}|n, \gamma_0, a = 0, p) \equiv p(z_{1:i}|i, \gamma_0, a = 0, p)$. Thus when $a = 0$, the EPPF becomes independent of the sample size, which is a well-known property for the Chinese restaurant process.

Corollary 4 (Sequential Construction). *Since $p(z_{i+1}|z_{1:i}, n, \gamma_0, a, p) = \frac{p(z_{1:i+1}|n, \gamma_0, a, p)}{p(z_{1:i}|n, \gamma_0, a, p)}$, conditioning on the sample size n , the sequential prediction rule of the generalized Chinese restaurant sampling formula $\mathbf{z}|n \sim \text{gCRSF}(n, \gamma_0, a, p)$ can be expressed as*

$$P(z_{i+1} = k | z_{1:i}, n, \gamma_0, a, p) = \begin{cases} (n_{k, (i)} - a) \frac{R_{n, \gamma_0, a, p}(i+1, l_{(i)})}{R_{n, \gamma_0, a, p}(i, l_{(i)})}, & \text{for } k = 1, \dots, l_{(i)}; \\ \gamma_0 p^{-a} \frac{R_{n, \gamma_0, a, p}(i+1, l_{(i)} + 1)}{R_{n, \gamma_0, a, p}(i, l_{(i)})}, & \text{if } k = l_{(i)} + 1; \end{cases} \quad (18)$$

where $i = 1, \dots, n - 1$.

With this sequential prediction rule, similar to an EPPF of Π , we can construct Π_{i+1} from Π_i in a sample of size n by assigning element $(i + 1)$ to $A_{z_{i+1}}$. When $a = 0$, we have

$$\frac{R_{n, \gamma_0, a, p}(i + 1, l_{(i)})}{R_{n, \gamma_0, a, p}(i, l_{(i)})} = \frac{R_{n, \gamma_0, a, p}(i + 1, l_{(i)} + 1)}{R_{n, \gamma_0, a, p}(i, l_{(i)})} = \frac{\Gamma(i + \gamma_0)}{\Gamma(i + 1 + \gamma_0)} = \frac{1}{i + \gamma_0},$$

and this sequential prediction rule becomes the same as that of a Chinese restaurant process with concentration parameter γ_0 .

Corollary 5. *The distribution of the number of clusters in $z_{1:i}$ in a sample of size n can be expressed as*

$$\begin{aligned} p(l_{(i)}|n, \gamma_0, a, p) &= p(l_{(i)}|i, \gamma_0, a, p) \frac{\sum_{\ell=0}^i \gamma_0^\ell p^{-a\ell} S_a(i, \ell)}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)} R_{n, \gamma_0, a, p}(i, l_{(i)}), \\ &= \frac{\gamma_0^{l_{(i)}} p^{-a l_{(i)}} S_a(i, l_{(i)}) R_{n, \gamma_0, a, p}(i, l_{(i)})}{\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell)}. \end{aligned} \quad (19)$$

This can be directly derived using (16) and the relationship between the EPPF and the distribution of the number of clusters. From this PMF, we obtain a useful identity

$$\sum_{\ell=0}^n \gamma_0^\ell p^{-a\ell} S_a(n, \ell) = \gamma_0 p^{-a} R_{n, \gamma_0, a, p}(1, 1),$$

which could be used to calculate the PMF of the generalized NB distribution in (8) and the EPPF in (13) without the need to compute the generalized Stirling numbers $S_a(n, l)$.

Corollary 6. *Given the model parameters γ_0 , a and p , the probability for two elements uniformly at random selected from a random sample of size n to be in two different groups can be expressed as*

$$P(z_1 \neq z_2 | n, \gamma_0, a, p) = \frac{\gamma_0 p^{-a} R_{n, \gamma_0, a, p}(2, 2)}{R_{n, \gamma_0, a, p}(1, 1)} = \left[1 + \frac{1-a}{\gamma_0 p^{-a}} \frac{R_{n, \gamma_0, a, p}(2, 1)}{R_{n, \gamma_0, a, p}(2, 2)} \right]^{-1}. \quad (20)$$

When $a = 0$, for $n \geq 2$, we have

$$P(z_1 \neq z_2 | n, \gamma_0, a = 0, p) \equiv \frac{\gamma_0}{1 + \gamma_0}.$$

Proof. We directly obtain (20) by setting $i = 1$ and $z_{i+1} = 2$ in (18) and using the recursive definition of $R_{n, \gamma_0, a, p}(1, 1)$ in (17). \square

Corollary 7 (Simpson's Index of Diversity). *Given the model parameters $\theta = \{\gamma_0, a, p\}$, the probability for two individuals uniformly at random selected from a random sample, whose size follows $n \sim \text{gNB}(\gamma_0, a, p)$ and is larger than two, to be in two different groups can be expressed as*

$$\begin{aligned} S_\theta &:= P(z_1 \neq z_2 | \gamma_0, a, p) = \sum_{n=2}^{\infty} P(z_1 \neq z_2 | n, \gamma_0, a, p) \frac{\text{gNB}(n; \gamma_0, a, p)}{1 - \text{gNB}(0; \gamma_0, a, p) - \text{gNB}(1; \gamma_0, a, p)} \\ &= \frac{\gamma_0^2 p^{-2a} e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}}}{1 - e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}} - \gamma_0 p^{1-a} e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}}} \sum_{n=2}^{\infty} \frac{p^n}{n!} R_{n, \gamma_0, a, p}(2, 2). \end{aligned} \quad (21)$$

When $a = 0$, we have

$$P(z_1 \neq z_2 | \gamma_0, a = 0, p) \equiv \frac{\gamma_0}{1 + \gamma_0}.$$

Under this construction, given a random species sample (z_1, \dots, z_n) , with a prior distribution on $\boldsymbol{\theta}$ as $p_{\Theta}(\boldsymbol{\theta})$, the posterior mean of Simpson's index of diversity is expressed as

$$S = \int S_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | z_1, \dots, z_n) d\boldsymbol{\theta}, \quad (22)$$

where

$$p(\boldsymbol{\theta} | z_1, \dots, z_n) = \frac{p(z_1, \dots, z_n, n | \boldsymbol{\theta}) p_{\Theta}(\boldsymbol{\theta})}{\int p(z_1, \dots, z_n, n | \boldsymbol{\theta}) p_{\Theta}(\boldsymbol{\theta}) d\boldsymbol{\theta}}.$$

In the next section we show how to perform MCMC estimation for the model from which we will derive the posterior value for Simpson's index of diversity.

5 Illustrations

Species abundance data of a sample is usually represented with a set of frequency counts $M = \{m_1, m_2, \dots\}$, where m_i denotes the number of species that have been observed i times in the sample. This data can also be converted into a sequence of group indices $\mathbf{z} = (z_1, \dots, z_n)$ or a group-size vector (n_1, \dots, n_l) , where n_k is the number of individuals in group k , $n = \sum_i i m_i = \sum_{k=1}^l n_k$ is the size of the sample and $l = \sum_i m_i$ is the number of distinct groups in the sample. For example, we may represent $M = \{m_1 = 2, m_2 = 1, m_3 = 2\}$ as $\mathbf{z} = (1, 2, 3, 3, 4, 4, 4, 5, 5, 5)$ or $(n_1, \dots, n_5) = (1, 1, 2, 3, 3)$. For a sample of species frequency counts, we use (12) as the likelihood for the model parameters $\boldsymbol{\theta} = \{\gamma_0, a, p\}$. With appropriate priors imposed on $\boldsymbol{\theta}$, we use MCMC to obtain posterior samples $\boldsymbol{\theta}^{(j)} = \{\gamma_0^{(j)}, a^{(j)}, p^{(j)}\}$ and then calculate $S_{\boldsymbol{\theta}^{(j)}}$. The details of MCMC update equations are provided in the Appendix.

5.1 Estimation of T-cell Receptor Diversity

An important characteristic of the immune system is the diversity of T-cell receptors (TCRs) (Nikolich-Zugich et al., 2004, Ferreira et al., 2009). As the number of distinct TCRs might be extremely high in the body, one usually investigates TCR diversity by collecting a sample of T-cells and determining the number of distinct TCR sequences and their respective abundances (counts) in that sample. For example, a Bayesian semiparametric approach is proposed in Guindani et al. (2014) to estimate TCR diversity of regulatory, Treg, and conventional T-cells, Tconv, across samples of two healthy and three diabetic mice; the TCR diversity there is defined as the number of distinct TCR sequences in a sample, including k'

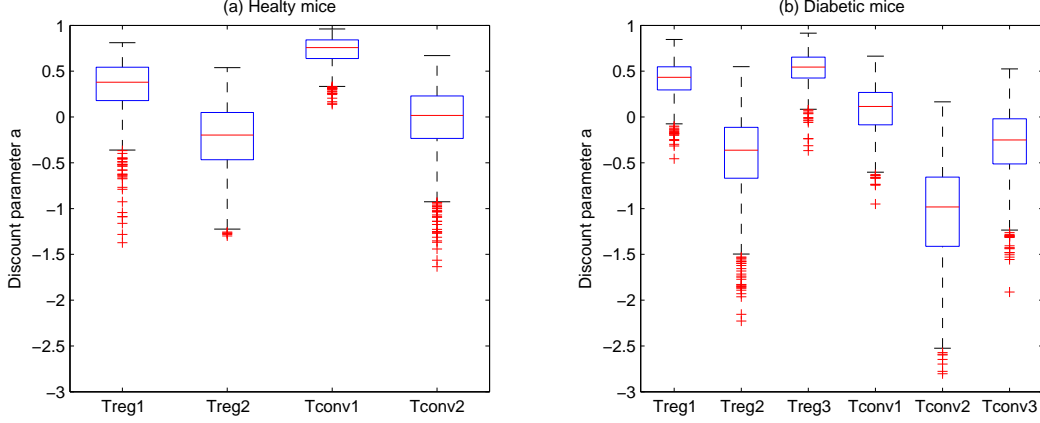


Figure 2: Box plots of $\{a^{(j)}\}_{j=1,N}$, the posterior MCMC samples of the discount parameter a , for regulatory, Treg, and conventional T-cells, Tconv, across various samples of (a) two healthy and (b) three diabetic mice.

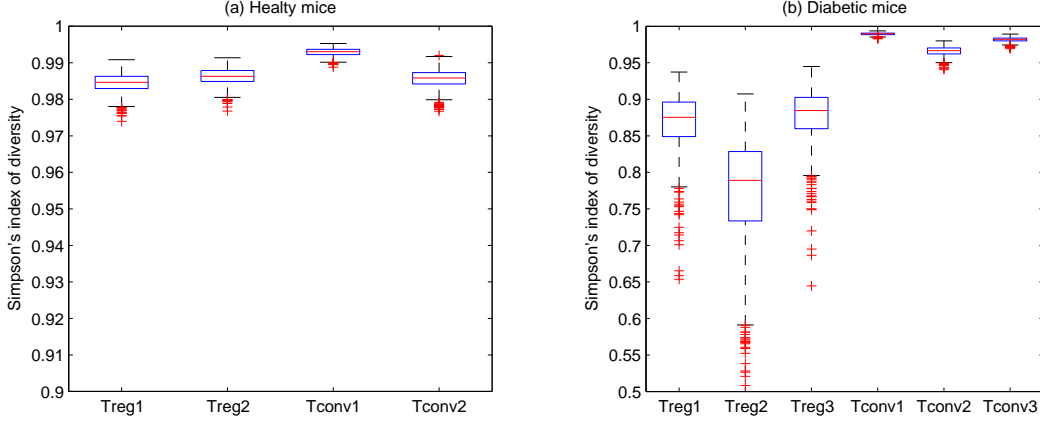


Figure 3: Box plots of $\{S_{\theta^{(j)}}\}_{j=1,N}$, the posterior MCMC samples of Simpson's index of diversity, for regulatory, Treg, and conventional T-cells, Tconv, across various samples of (a) two healthy and (b) three diabetic mice.

observed distinct TCR sequences and k_0 unobserved ones due to censoring of zero counts. In this paper, we estimate TCR diversity by calculating Simpson's index of diversity given a sample of species frequency counts.

Considering the same TCR species abundance frequency count dataset used in Ferreira et al. (2009) and presented in Table 2 of Guindani et al. (2014), we compare Simpson's index of diversity of the TCRs of Treg and Tconv across samples of two healthy and three diabetic mice. For example, for Treg, we have $M = \{40, 5, 5, 2, 3\}$ with $i \in \{1, 2, 3, 4, 5\}$ for the sample from healthy mouse 1, and we have $M = \{8, 1, 2, 1, 1, 1\}$ with $i \in \{1, 2, 3, 5, 36, 40\}$ for the sample from diabetic mouse 1. For each sample of T-cells, we consider 2000 MCMC iterations and collect the last 1000 MCMC samples $\{\theta^{(j)}\}_{1,1000}$.

Figure 2 shows the box plots of the MCMC posterior samples of the discount parameter a in various samples of regulatory and conventional T-cells for the healthy and diabetic mice. We find no clear associations between the posteriors of a and whether the mice are healthy or diabetic or whether the T-cells are regulatory or conventional.

As shown in Figure 3, using the samples for the diabetic mice, the estimated Simpson’s indices of diversity of the TCRs for regulatory T-cells are considerably lower than those for conventional T-cells; whereas for the healthy mice, no clear differences on TCR diversity are found. Comparing Figures 2 and 3, one may also not find clear relationships between the estimated values of a and the estimated Simpson’s indices of diversity, which suggests that for the generalized negative binomial process, the discount parameter a alone may not be a good indicator for species evenness measured by Simpson’s index of diversity. Guindani et al. (2014) showed that diabetic mice tended to have a smaller number of distinct TCRs in a sample of regulatory T-cells than in a sample of conventional T-cells. Our comparison of Simpson’s indices of diversity, which measure species evenness and hence complementary to the comparison of species richness studied in Guindani et al. (2014), provides additional evidence to suggest that for diabetic mice, the TCR diversity of regulatory T-cells is lower than that of conventional T-cells.

5.2 Genomic Data Analysis

An important research topic in genomics is the analysis of expressed sequence tag (EST) data, which arise by sequencing complementary DNA (cDNA) libraries consisting of millions of genes. The number of ESTs from a particular gene indicates the expression level of that gene. It is typical that only a small portion of the cDNA is sequenced in a sample due to cost constraints, and one needs to rely on this sample to estimate population properties. We consider a tomato flower EST dataset, previously analyzed in Mao and Lindsay (2002) and Lijoi et al. (2007a), that consists of 2586 ESTs from 1825 genes as $M = \{1434, 253, 71, 33, 11, 6, 2, 3, 1, 2, 2, 1, 1, 1, 2, 1, 1\}$ for $i \in \{1, \dots, 14\} \cup \{16, 23, 27\}$. We convert $\{m_i\}_i$ into (z_1, \dots, z_{2586}) . To evaluate the accuracy of the proposed nonparametric Bayesian estimator in (21), we consider this relatively large sample as the population and treat $\hat{S} = 0.9993$, a sample estimate with (2), as the “true” Simpson’s index of diversity for the population.

We randomly select an EST sample of size $n = 50$ from (z_1, \dots, z_{2586}) to estimate the Simpson’s index of diversity of the population. For each selected EST sample, we use MCMC to obtain posterior samples $\boldsymbol{\theta}^{(j)} = \{\gamma_0^{(j)}, a^{(j)}, p^{(j)}\}$ and then calculate $S_{\boldsymbol{\theta}^{(j)}}$; we consider 2000 MCMC iterations and collect one sample in every five iterations in the last 1000 MCMC

Table 1: Simulation study based on 100 expressed sequence tag (EST) samples of size 50 uniformly at random selected from a population of 2586 ESTs from 1825 distinct genes, with various settings of the discount parameter a . A sample estimate of 0.9993 using all the 2586 ESTs is considered as the “true” Simpson’s index of diversity for the population.

Parameter Setting	Mean Bias ($\times 10^{-3}$)	Median Bias ($\times 10^{-3}$)	50% Coverage	95% Coverage
$a = -1$	10.37	10.60	0%	0%
$a = 0$	3.05	3.31	0%	0%
$a = 0.5$	1.07	1.40	18%	85%
$a < 0$	3.51	3.78	0%	0%
$0 \leq a < 1$	0.48	1.11	62%	98%
$a < 1$	0.41	1.09	69%	99%

iterations, leading to $N = 200$ total samples $\{\theta^{(j)}\}_{1,200}$; we find from the collected MCMC samples the mean, median, 50 percentile range and 95 percentile range of $\{S_{\theta^{(j)}}\}$, and compare these values against 0.9993. We repeat the same procedure 100 times and find the averages among these 100 times of the absolute distances from the mean and median to 0.9993, and the probabilities for 0.9993 to be covered by the 50 and 95 percentile ranges.

We summarize the results in Table 1, where we fix a to be -1 , 0 or 0.5 , or let a be inferred for each EST sample and restrict it to be $a < 0$, $0 \leq a < 1$ or $a < 1$. It is clear that allowing a to be freely adjusted within $(-\infty, 1)$ leads to a more accurate estimation of Simpson’s index of diversity using a sample of the population, demonstrating the effectiveness of the generalized negative binomial process on the analysis of EST sequence counts. Similar simulation results are observed on the TCR sequence count dataset studied in Section 5.1.

In conclusion, we have introduced a sample size dependent species model, which allows flexible modeling of species abundance frequency count data. We gain this flexibility with a simple model and consequently posterior inference via MCMC is also simple. The paper provides a framework to jointly model a single random count and its exchangeable random partition. It is natural to extend the same framework to mixture modeling, where the usual task is to partition a set of data points into exchangeable clusters, where both the number and sizes of clusters are unknown and need to be inferred. The techniques developed here to model a random count vector also serve as the foundation for Zhou et al. (2014) to construct a family of nonparametric Bayesian priors for infinite random count matrices, and for Zhou (2014) to define a prior distribution that describes the random partition of a count vector into a latent random count matrix.

References

- D. Aldous. Exchangeability and related topics. In *Ecole d'Ete de Probabilities de Saint-Flour XIII*, pages 1–198. Springer, 1983.
- C. Antoniak. Mixtures of Dirichlet processes with applications to bayesian nonparametric problems. *The Annals of Statistics*, (2):1152–1174, 1974.
- E. Barrios, A. Lijoi, L. E. Nieto-Barajas, and I. Pruenster. Modeling with normalized random measure mixture models. *Carlo Alberto Notebooks*, No. 276, 2012.
- A. Brix. Generalized gamma measures and shot-noise Cox processes. *Advances in Applied Probability*, 1999.
- J. Bunge and M. Fitzpatrick. Estimating the number of species: a review. *Journal of the American Statistical Association*, 1993.
- J. Bunge, A. Willis, and F. Walsh. Estimating the number of species in microbial diversity studies. *Annual Review of Statistics and Its Application*, 2014.
- A. Cerquetti. Bayesian nonparametric estimation of Simpson’s evenness index under alpha-Gibbs priors. *arXiv:1203.1666*, 2012.
- A. Chao. Species richness estimation. *Encyclopedia of statistical sciences*, 12:7907–7916, 2005.
- C. A Charalambides. *Combinatorial methods in discrete distributions*. Wiley, 2005.
- S. Engen. On species frequency models. *Biometrika*, 1974.
- W. J. Ewens. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, 1972.
- S. Favaro and Y. W. Teh. MCMC for normalized random measure mixture models. *to appear in Statistical Science*, 2013.
- S. Favaro, A. Lijoi, R. H. Mena, and I. Prünster. Bayesian non-parametric inference for species variety with a two-parameter Poisson–Dirichlet process prior. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2009.
- S. Favaro, A. Lijoi, and I. Pruenster. Conditional formulae for gibbs-type exchangeable random partitions. *Annals of Applied Probability*, 2013.
- T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1973.
- T. S. Ferguson and M. J. Klass. A representation of independent increment processes without gaussian components. *Annals of Mathematical Statistics*, 1972.
- C. Ferreira, Y. Singh, A. L. Furmanski, F. S. Wong, O. A. Garden, and J. Dyson. Non-obese diabetic mice select a low-diversity repertoire of natural regulatory t cells. *Proceedings of the National Academy of Sciences*, 2009.
- R. A. Fisher, A. Steven Corbet, and C. B. Williams. The relation between the number of species and the number of individuals in a random sample of an animal population. *Journal of Animal Ecology*, 1943.
- H. U. Gerber. From the generalized gamma to the generalized negative binomial distribution. *Insurance: mathematics and economics*, 1992.

- A. Gnedin and J. Pitman. Exchangeable Gibbs partitions and Stirling triangles. *Journal of Mathematical Sciences*, 2006.
- A. Gnedin, C. Haulk, and J. Pitman. Characterizations of exchangeable partitions and random discrete distributions by deletion properties. In N.H. Bingham and C.M. Goldie, editors, *Probability and Mathematical Genetics: Papers in Honour of Sir John Kingman*. 2009.
- J. E. Griffin and S. G. Walker. Posterior simulation of normalized random measure mixtures. *Journal of Computational and Graphical Statistics*, 2011.
- M. Guindani, N. Sepúlveda, C. D. Paulino, and P. Müller. A bayesian semiparametric approach for the differential analysis of sequence counts data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 2014.
- M. O. Hill. Diversity and evenness: a unifying notation and its consequences. *Ecology*, 1973.
- J. B. Hughes, J. J. Hellmann, T. H. Ricketts, and B. J. M. Bohannan. Counting the uncountable: statistical approaches to estimating microbial diversity. *Applied and Environmental Microbiology*, 2001.
- H. Ishwaran and L. F. James. Gibbs sampling methods for stick-breaking priors. *JASA*, 2001.
- L. F. James, A. Lijoi, and I. Prünster. Posterior analysis for normalized random measures with independent increments. *Scandinavian Journal of Statistics*, 2009.
- J. F. C. Kingman. Completely random measures. *Pacific Journal of Mathematics*, 1967.
- J. F. C. Kingman. Random partitions in population genetics. *Proceedings of the Royal Society of London. A.*, 1978a.
- J. F. C. Kingman. The representation of partition structures. *Journal of the London Mathematical Society*, 1978b.
- J. F. C. Kingman. *Poisson Processes*. Oxford University Press, 1993.
- A. Lijoi and I. Prünster. Models beyond the Dirichlet process. In N. L. Hjort, C. Holmes, P. Müller, and S. G. Walker, editors, *Bayesian nonparametrics*. Cambridge University Press, 2010.
- A. Lijoi, R. H. Mena, and I. Prünster. Hierarchical mixture modeling with normalized inverse-Gaussian priors. *Journal of the American Statistical Association*, 2005.
- A. Lijoi, R. H. Mena, and I. Prünster. Bayesian nonparametric estimation of the probability of discovering new species. *Biometrika*, 2007a.
- A. Lijoi, R. H. Mena, and I. Prünster. Controlling the reinforcement in Bayesian non-parametric mixture models. *Journal of the Royal Statistical Society: Series B*, 2007b.
- A. Lijoi, I. Prünster, and S. G. Walker. Bayesian nonparametric estimators derived from conditional Gibbs structures. *The Annals of Applied Probability*, 2008.
- A. E. Magurran. *Measuring biological diversity*. Taylor & Francis, 2004.
- C. X. Mao and B. G. Lindsay. A Poisson model for the coverage problem with a genomic application. *Biometrika*, 2002.
- P. Müller and R. Mitra. Bayesian nonparametric inference – why and how. *Bayesian Analysis*, 2013.

- P. Müller and F. A. Quintana. Nonparametric Bayesian data analysis. *Statistical Science*, 2004.
- J. Nikolich-Zugich, M. K. Slifka, and I. Messaoudi. The many important facets of T-cell repertoire diversity. *Nature Reviews Immunology*, 2004.
- M. Perman, J. Pitman, and M. Yor. Size-biased sampling of poisson point processes and excursions. *Probability Theory and Related Fields*, 1992.
- J. Pitman. Exchangeable and partially exchangeable random partitions. *Probability Theory and Related Fields*, 1995.
- J. Pitman. Some developments of the Blackwell-Macqueen urn scheme. In *Statistics, Probability and Game Theory; Papers in honor of David Blackwell*, 1996.
- J. Pitman. Poisson-Kingman partitions. *Lecture Notes-Monograph Series*, pages 1–34, 2003.
- J. Pitman. *Combinatorial stochastic processes*. Lecture Notes in Mathematics. Springer-Verlag, 2006.
- J. Pitman and M. Yor. The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *The Annals of Probability*, 1997.
- M. H. Quenouille. A relation between the logarithmic, Poisson, and negative binomial series. *Biometrics*, 1949.
- E. Regazzini, A. Lijoi, and I. Prünster. Distributional results for means of normalized random measures with independent increments. *Annals of Statistics*, 2003.
- C. Ritter and M. A. Tanner. Facilitating the Gibbs sampler: the Gibbs stopper and the griddy-Gibbs sampler. *Journal of the American Statistical Association*, 1992.
- A. K. Shaw, A. L. Halpern, K. Beeson, B. Tran, J. C. Venter, and J. B. H. Martiny. It’s all relative: ranking the diversity of aquatic bacterial communities. *Environmental microbiology*, pages 2200–2210, 2008.
- E. H. Simpson. Measurement of diversity. *Nature*, 1949.
- G. E. Willmot. A remark on the poisson-pascal and some other contagious distributions. *Statistics & probability letters*, 1988.
- M. Zhou. Beta-negative binomial process and exchangeable random partitions for mixed-membership modeling. *To appear in NIPS*, 2014.
- M. Zhou and L. Carin. Negative binomial process count and mixture modeling. *To appear in IEEE Trans. Pattern Analysis and Machine Intelligence*, 2013.
- M. Zhou, L. Hannah, D. Dunson, and L. Carin. Beta-negative binomial process and Poisson factor analysis. In *AISTATS*, 2012.
- M. Zhou, O.-H. Madrid-Padilla, and J. G. Scott. Priors for random count matrices derived from a family of negative binomial processes. *arXiv:1404.3331v2*, 2014.

A Proof for Theorem 1

Proof. Let us consider the process X_G , conditional on G , given by

$$X_G(A) = \sum_k n_k \mathbf{1}(\omega_k \in A).$$

Now it is easy to see that

$$\mathbb{E}[\exp\{-\phi X_G(A)\} | G] = \exp\{-G(A)(1 - e^{-\phi})\},$$

and using the well known result for homogeneous Lévy processes, we have

$$\mathbb{E}[\exp\{-\lambda G(A)\}] = \exp\left\{-G_0(A) \int_0^\infty [1 - e^{-\lambda s}] \rho(ds)\right\}. \quad (23)$$

Now, the key observation is the following identity:

$$1 - e^{-(1-e^{-\phi})s} = 1 - e^{-s} \sum_{j=0}^\infty \frac{s^j}{j!} e^{-\phi j} = (1 - e^{-s}) - e^{-s} \sum_{j=1}^\infty \frac{s^j}{j!} e^{-\phi j}.$$

Let us put this to one side for now and consider the model for \tilde{X} given by

$$\tilde{X}(A) = \sum_{k=1}^l n_k \mathbf{1}(\omega_k \in A)$$

with $l \sim \text{Po}(\gamma G_0(\Omega))$ for some non-negative γ and independently $P(n_k = j) = \pi_j$ for some $\pi_j \leq 1$ and $j \in \{1, 2, \dots\}$. Now given l , we have

$$\mathbb{E}[\exp\{-\phi \tilde{X}(A)\} | l] = \prod_{k=1}^l \mathbb{E}[\exp\{-\phi n_k \mathbf{1}(\omega_k \in A)\}]$$

and each of these expectations is given by

$$\psi = \sum_{j=1}^\infty e^{-\phi j} \pi_j.$$

Thus

$$\mathbb{E}[\exp\{-\phi \tilde{X}(A)\}] = \exp\{-\gamma G_0(A) (1 - \psi)\}$$

which is given by

$$\exp \left\{ -\gamma G_0(A) \left[1 - \sum_{j=1}^{\infty} e^{-\phi j} \pi_j \right] \right\}. \quad (24)$$

Comparing (23) and (24) we see that we have a match when

$$\gamma = \int_0^{\infty} (1 - e^{-s}) \rho(ds)$$

and

$$\pi_j = \frac{\int_0^{\infty} s^j e^{-s} \rho(ds)}{j! \gamma},$$

and note that it is easy to verify that

$$\sum_{j=1}^{\infty} \pi_j = 1.$$

□

B Proof for Corollary 3

This follows directly from Bayes' rule, since $p(z_i | \mathbf{z}^{-i}, n, \gamma_0, \rho) = \frac{p(z_i, \mathbf{z}^{-i}, n | \gamma_0, \rho)}{p(\mathbf{z}^{-i}, n | \gamma_0, \rho)}$, where

$$p(z_i, \mathbf{z}^{-i}, n | \gamma_0, \rho) = n^{-1} p(\mathbf{z}^{-i}, n-1 | \gamma_0, \rho) \left[\gamma_0 \int_0^{\infty} s e^{-s} \rho(ds) \mathbf{1}(z_i = l^{-i} + 1) + \sum_{k=1}^{l^{-i}} \frac{\int_0^{\infty} s^{n_k^{-i}+1} e^{-s} \rho(ds)}{\int_0^{\infty} s^{n_k^{-i}} e^{-s} \rho(ds)} \mathbf{1}(z_i = k) \right].$$

Marginalizing out the z_i from $p(z_i, \mathbf{z}^{-i}, n | \gamma_0, \rho)$ we have

$$p(\mathbf{z}^{-i}, n | \gamma_0, \rho) = n^{-1} p(\mathbf{z}^{-i}, n-1 | \gamma_0, \rho) \left[\gamma_0 \int_0^{\infty} s e^{-s} \rho(ds) + \sum_{k=1}^{l^{-i}} \frac{\int_0^{\infty} s^{n_k^{-i}+1} e^{-s} \rho(ds)}{\int_0^{\infty} s^{n_k^{-i}} e^{-s} \rho(ds)} \right].$$

C Derivations for the GNB

Marginalizing out λ from $[n | \lambda] \sim \text{Po}(\lambda)$ with $\lambda \sim \text{gGamma}[\gamma_0, a, p/(1-p)]$, leads to a generalized NB distribution; $n \sim \text{gNB}(\gamma_0, a, p)$, with shape parameter γ_0 , discount parameter $a < 1$, and probability parameter p . The probability generating function (PGF) is given by

$$\mathbb{E}[t^n] = \mathbb{E}[\mathbb{E}[t^n | \lambda]] = \exp \left\{ -\frac{\gamma_0 [(1-pt)^a - (1-p)^a]}{ap^a} \right\},$$

the mean value is $\gamma_0 [p/(1-p)]^{1-a}$ and the variance is $\gamma_0 [p/(1-p)]^{1-a} (1-ap)/(1-p)$. The PGF was originally presented in Willmot (1988) and Gerber (1992). With the PGF written as

$$\begin{aligned}\mathbb{E}(t^n) &= \exp \left\{ \gamma_0 \frac{(1-p)^a}{ap^a} \right\} \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{-\gamma_0(1-pt)^a}{ap^a} \right)^k \\ &= \exp \left\{ \gamma_0 \frac{(1-p)^a}{ap^a} \right\} \sum_{k=0}^{\infty} \frac{1}{k!} \left(\frac{-\gamma_0}{ap^a} \right)^k \sum_{j=0}^{\infty} \binom{ak}{j} (-pt)^j,\end{aligned}$$

we can derive the PMF as

$$p_N(n|\gamma_0, a, p) = \frac{p^n}{n!} e^{\gamma_0 \frac{(1-p)^a}{ap^a}} \sum_{k=0}^{\infty} \frac{1}{k!} \left(-\frac{\gamma_0}{ap^a} \right)^k \frac{\Gamma(n-ak)}{\Gamma(-ak)}, \quad n = 0, 1, \dots \quad (25)$$

We can also generate $n \sim \text{gNB}(\gamma_0, a, p)$ from a compound Poisson distribution, as $n = \sum_{k=1}^l n_k$, with the (n_k) independent from $\text{TNB}(a, p)$, and $l \sim \text{Po}\left(\frac{\gamma_0(1-(1-p)^a)}{ap^a}\right)$, where $\text{TNB}(a, p)$ denotes a truncated NB distribution, with PGF $\mathbb{E}[t^u] = \frac{1-(1-pt)^a}{1-(1-p)^a}$ and PMF

$$p_U(u|a, p) = \frac{\Gamma(u-a)}{u!\Gamma(-a)} \frac{p^u(1-p)^{-a}}{1-(1-p)^{-a}}, \quad u = 1, 2, \dots \quad (26)$$

Note that as $a \rightarrow 0$, $u \sim \text{TNB}(a, p)$ becomes a logarithmic distribution (Quenouille, 1949) with PMF $p_U(u|p) = \frac{-1}{\ln(1-p)} \frac{p^u}{u}$ and $n \sim \text{gNB}(\gamma_0, a, p)$ becomes a NB distribution; $n \sim \text{NB}(\gamma_0, p)$. The truncated NB distribution with $0 < a < 1$ is the extended NB distribution introduced in Engen (1974).

Here we provide a useful identity which we will be used later in this section. Denote by \sum_* as the summation over all sets of positive integers (n_1, \dots, n_l) with $\sum_{k=1}^l n_k = n$. We call $n \sim \text{SumTNB}(l, a, p)$ as a sum-truncated NB distributed random variable that can be generated via $n = \sum_{k=1}^l n_k$, $n_k \sim \text{TNB}(a, p)$. Using both (26) and

$$\left[\frac{1-(1-pt)^a}{1-(1-p)^a} \right]^l = \frac{\sum_{k=0}^l \binom{l}{k} (-1)^k \sum_{j=0}^{\infty} \binom{ak}{j} (-pt)^j}{[1-(1-p)^a]^l},$$

we may express the PMF of the sum-truncated NB distribution as

$$p_N(n|l, a, p) = \sum_* \prod_{k=1}^l \frac{\Gamma(n_k-a)}{n_k!\Gamma(-a)} \frac{p^{n_k}(1-p)^{-a}}{1-(1-p)^{-a}} = \frac{p^n}{[1-(1-p)^a]^l} \sum_{k=0}^l (-1)^k \binom{l}{k} \frac{\Gamma(n-ak)}{n!\Gamma(-ak)},$$

leading to the identity shown in (9).

The EPPF is the ECPF in (12) divided by the marginal distribution of n in (25), given

by

$$p(\mathbf{z}|n, \gamma_0, a, p) = p_n(z_1, \dots, z_n|n) = \frac{e^{-\frac{\gamma_0}{ap^a}}}{\sum_{k=0}^{\infty} \frac{1}{k!} \left(-\frac{\gamma_0}{ap^a}\right)^k \frac{\Gamma(n-ak)}{\Gamma(-ak)}} \gamma_0^l p^{-al} \prod_{k=1}^l \frac{\Gamma(n_k - a)}{\Gamma(1 - a)}. \quad (27)$$

Using the EPPF in (13) and the identity in (9), the conditional distribution of the number of clusters l in a sample of size n can be expressed as

$$p_L(l|n, \gamma_0, a, p) = \frac{1}{l!} \sum_{*} \frac{n!}{\prod_{k=1}^l n_k!} p(\mathbf{z}|n, \gamma_0, a, p) = \frac{\gamma_0^l p^{-al} S_a(n, l)}{e^{\frac{\gamma_0}{ap^a}} \sum_{k=0}^{\infty} \frac{1}{k!} \left(-\frac{\gamma_0}{ap^a}\right)^k \frac{\Gamma(n-ak)}{\Gamma(-ak)}}, \quad (28)$$

which, since $\sum_{l=0}^n p_L(l|n, \gamma_0, a, p) = 1$, further leads to identity

$$e^{\frac{\gamma_0}{ap^a}} \sum_{k=0}^{\infty} \frac{1}{k!} \left(-\frac{\gamma_0}{ap^a}\right)^k \frac{\Gamma(n-ak)}{\Gamma(-ak)} = \sum_{l=0}^n \gamma_0^l p^{-al} S_a(n, l).$$

Applying this identity on (25), (27) and (28) lead to (8), (13) and (14).

D MCMC Inference

For the gNBP, the ECPF in (12) defines a fully factorized likelihood for γ_0 , a and p . With a gamma prior $\text{Gamma}(e_0, 1/f_0)$ placed on γ_0 , we have

$$(\gamma_0|-) \sim \text{Gamma}\left(e_0 + l, \frac{1}{f_0 + \frac{1-(1-p)^a}{ap^a}}\right). \quad (29)$$

As $a \rightarrow 0$, we have $(\gamma_0|-) \sim \text{Gamma}\left(e_0 + l, \frac{1}{f_0 - \ln(1-p)}\right)$. This paper sets $e_0 = f_0 = 0.01$.

Since $a < 1$, we have $\tilde{a} = \frac{1}{1+(1-a)} \in (0, 1)$. With a uniform prior placed on \tilde{a} in $(0, 1)$ and the likelihood of gNBP in (12), we use the griddy-Gibbs sampler (Ritter and Tanner, 1992) to sample a from a discrete distribution

$$P(a|-) \propto e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}} p^{-al} \prod_{k=1}^l \frac{\Gamma(n_k - a)}{\Gamma(1 - a)} \quad (30)$$

over a grid of points $\frac{1}{1+(1-a)} = 0.0001, 0.0002, \dots, 0.9999$.

We place a uniform prior on p in $(0, 1)$. When $a \rightarrow 0$, the likelihood of the gNBP in (12) becomes proportional to $p^m (1-p)^{\gamma_0}$, thus we have $(p|-) \sim \text{Beta}(1+n, 1+\gamma_0)$. When $a \neq 0$, we use the griddy-Gibbs sampler to sample p from a discrete distribution

$$P(p|-) \propto e^{-\gamma_0 \frac{1-(1-p)^a}{ap^a}} p^{n-al} \quad (31)$$

over a grid of points $p = 0.001, 0.002, \dots, 0.999$.